

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

EV205822953

**Community Mining Based On Core Objects And
Affiliated Objects**

Inventor(s):
Ji-Rong Wen
Wen-Jun Zhou
Wei-Ying Ma
Hong-Jiang Zhang

ATTORNEY'S DOCKET NO. MS1-1617US

1 **TECHNICAL FIELD**

2 This invention relates to community mining, and particularly to community
3 mining based on core objects and affiliated objects.
4

5 **BACKGROUND**

6 Discovering related objects from a collection of objects is a very useful
7 capability, particularly when the collection of objects becomes very large. This
8 problem of discovering similar or related objects from a collection of objects is
9 also referred to as community mining. By mining communities of related objects
10 from a large collection of objects, groups of related objects are able to be
11 identified more quickly and easily than when using other methodologies (such as
12 manually scanning all of the objects in the collection).

13 By way of example, a large number of web pages exists on the Internet. It
14 would be useful to be able to group these web pages together into communities of
15 related web pages, allowing users to quickly and easily view these communities.
16 By way of another example, a database of papers written by researchers may be
17 available. It would be useful to be able to group these papers together into
18 communities of related papers, allowing users to quickly and easily view these
19 communities.

20 The accuracy of current community mining techniques, however, is lacking.
21 Thus, it would be beneficial to improve on the manner in which community
22 mining is performed.
23
24
25

SUMMARY

Community mining based on core objects and affiliated objects is described herein.

According to one aspect, a set of core objects for a community of objects are identified from a plurality of objects. The community is expanded, based on the set of core objects, to include a set of affiliated objects.

According to another aspect, a model of a community of objects is obtained by grouping a first collection of a plurality of objects into a center portion, and grouping a second collection of the plurality of objects into one or more concentric portions around the center portion. The groupings of the first and second collections of the objects are identified as the community of objects.

BRIEF DESCRIPTION OF THE DRAWINGS

The same numbers are used throughout the document to reference like components and/or features.

Fig. 1 illustrates an example environment in which the community mining described herein can be implemented.

Fig. 2 is a flowchart illustrating an example process for community mining.

Fig. 3 illustrates an example concentric ring model of a community mined from a collection of objects.

Fig. 4 is a flowchart illustrating an example process for finding core sets of objects.

Fig. 5 shows an example that illustrates the process of Fig. 4.

Fig. 6 is a flowchart illustrating an example process for expanding communities based on core sets of objects.

1 Fig. 7 is a flowchart illustrating an example process for performing core set
2 merging.

3 Fig. 8 illustrates an example of core set merging.

4 Fig. 9 is a flowchart illustrating an example process for performing
5 community merging.

6 Fig. 10 illustrates an example of community merging.

7 Fig. 11 illustrates a general computer environment.

8 9 **DETAILED DESCRIPTION**

10 Community mining based on core objects and affiliated objects is described
11 herein. Sets of core objects are identified from a collection of objects, and these
12 sets of core objects are used as a basis for expansion and identification of affiliated
13 objects. A set of core objects and its affiliated objects collectively represent a
14 community. The affiliated objects can further be ranked in accordance with how
15 related they are to the set of core objects in the community. In one aspect, a
16 concentric circle model of the community is defined, with the set of core objects
17 being in the center of the concentric circle model, and each concentric circle
18 surrounding the center including affiliated objects having the same rank.

19 Fig. 1 illustrates an example environment 100 in which the community
20 mining described herein can be implemented. An object collection description
21 102 is accessible to a community mining module 104. Once description 102 is
22 accessed, community mining module 104 discovers one or more communities
23 from the described collection of objects 102, and outputs the one or more
24 discovered communities 106.
25

1 Description 102 can be made accessible to module 104 in a variety of
2 different manners. For example, description 102 may be passed to module 104, or
3 module 104 may be made aware of the location of description 102 (e.g., in a
4 database) and retrieve description 102 from that location. Alternatively, the
5 objects of the collection themselves may be passed to or retrieved by module 104
6 and module 104 may generate the description.

7 Communities of related objects are automatically discovered by community
8 mining module 104. Communities can be mined for any of a wide variety of types
9 of objects by module 104. Examples of such types of objects include people,
10 documents, web pages, and so forth. The number of communities mined from a
11 collection of objects can vary based on the design of module 104 as well as the
12 particular objects in the collection, as discussed in more detail below.

13 The object collection description 102 describes the collection of objects and
14 the relationships (also referred to herein as links) between the objects. The
15 description 102 can be thought of as a graph with multiple nodes and lines
16 connecting some of the nodes. Each node of the graph represents an object in the
17 collection, and each line between two nodes in the graph represents a relationship
18 between the two nodes of the graph.

19 The exact nature of the link between two objects depends on the type of
20 objects. For example if the objects are people then the link between objects can be
21 a social relationship between the people (e.g., the two people are friends or
22 acquaintances, the two people are related to one another (e.g., part of the same
23 family by blood or by some legal means such as marriage or adoption), the two
24 people are co-workers, etc.). By way of another example, if the objects are
25 documents, then the link between objects can be a citation of one document in the

1 other. By way of yet another example, if the objects are web pages, then the link
2 between objects can be a hyperlink to one document being included in the other
3 document.

4 References are made herein to one object being linked to another object.
5 When a first object is linked to a second object, then that first object has a
6 relationship with the second object (e.g., if the objects are documents then the first
7 object may have a citation to the second object; if the objects are web pages then
8 the first object may have a hyperlink to the second object). Thus, for example, if a
9 particular document is cited by three other documents, then each of those other
10 three documents is linked to that particular document.

11 The links between objects can also be assigned weights, or mapped to a
12 numeric value in other words, to identify the difference among links. Various
13 mapping functions (e.g., a 0/1 function, a reciprocal function, etc.) could be
14 defined and used to represent such differences. The difference among links can be
15 the result of the difference among objects, or among the relationship definition
16 itself. For example, assume that document A is cited by an important document B
17 which is cited by numerous documents. Further assume that document A is also
18 cited by a not so important document C which is cited by no other documents.
19 The citation relationship, or the link, between A and B could be assigned a higher
20 value than that of the link between A and C. Another example is in a social
21 network. The marriage of two persons could be mapped to a higher value than the
22 co-worker relationship, thus representing the difference between close degree of
23 the human relationship.

24 It should also be noted that the link between two objects can be a direct link
25 or an indirect link. A direct link refers to the situation where no other objects

1 being situated in the path between the two objects. An indirect link refers to the
2 situation where one or more other objects are situated in the path between two
3 objects (e.g., if there's a direct link from object A to object B, and a direct link
4 from object B to object C, then there is also an indirect link from object A to
5 object C (with object B being situated in the path from object A to object C)).
6 References herein to links can refer to direct links and/or indirect links (which
7 links are to be used can vary by implementation as desired by the designer or user
8 of community mining module 104 and/or the generator of object collection
9 description 102).

10 The description 102 of the object collection and the relationships between
11 objects in the collection may be generated by another device or component and
12 passed to (or otherwise made available to) community mining module 104, or
13 alternatively may be generated by module 104. The manner in which the
14 relationships between objects in a collection are determined can vary based on the
15 type of objects in the collection. For example, if the objects are people then the
16 links between people can be determined based on information collected from the
17 people (e.g., via questionnaires, registration processes, publicly available
18 information, etc.), or by analyzing people's personal web pages to identify
19 references to other people's web pages; if the objects are web pages then the links
20 between web pages can be determined by searching each web page for hyperlinks
21 to other web pages; and if the objects are documents then the links between
22 documents can be determined by searching each document for citations to other
23 documents.

24 Community mining module 104 uses the description of the collection of
25 objects and the links between the objects to discover the communities within the

1 collection of objects. Community mining module 104 includes a core set
2 identification module 112, a community expansion module 114, an optional core
3 set merging module 116, and an optional community merging module 118.

4 Core set identification module 112 identifies groups or sets of core objects
5 for communities. The core objects for a particular community are collectively
6 referred to as a core object set and as the core object set for that particular
7 community. These core objects are typically objects that are linked to by large
8 numbers of other objects (e.g., documents that are frequently cited by other
9 documents, or web pages that are frequently hyperlinked to by other web pages).
10 The core objects serve as the core or center of a community.

11 Community expansion module 114 expands the communities with
12 additional objects, referred to as affiliated objects, based on the core objects. In
13 one implementation, each of the affiliated objects in a community is linked to one
14 or more of the core objects in the community.

15 Core set merging module 116 attempts to merge core sets based on the
16 similarity between the core sets. Situations can arise where two core sets are
17 identified in the collection of objects and, due to the particular links among the
18 objects, the two core sets may be very similar. If there is sufficient overlap or
19 similarity among objects in the core sets, then the two core sets are merged by core
20 set merging module 116, and the two communities having those two core sets
21 become a single community.

22 Community merging module 118 attempts to merge communities based on
23 the similarity between the objects in the communities. Situations can arise where
24 two communities are discovered in the collection of objects and, due to the
25 particular links among the objects, the two communities may be very similar. If

1 there is sufficient overlap or similarity among objects in the community (core set
2 objects as well as affiliated objects), then the two communities are merged by
3 community merging module 118, and the two communities become a single
4 community.

5 It should be noted that community mining module 104 can be implemented
6 in any of a variety of manners. For example, module 104 may be implemented on
7 a general purpose computing device, or alternatively on a specialized computing
8 device (e.g., specialized for community mining). Additionally, it is to be
9 appreciated that the different modules 112, 114, 116, and 118 may all be
10 implemented on the same device or alternatively may be distributed across
11 multiple devices, and furthermore that the functionality of the individual modules
12 112, 114, 116, and/or 118 may also be distributed across multiple devices.

13 In certain embodiments, the mined communities 106 output by community
14 mining module 104 are defined as a four-tuple $\langle C, A, F, Va \rangle$, where C represents
15 the set of core objects, A represents the set of affiliated objects, F represents the
16 affiliation definition function measuring two objects i and j (which will return a
17 positive value if i is affiliated by j , such as a value of 1 if j has a direct link to i and
18 a value of 0 otherwise, or a function defined under a complex weighted graph that
19 if there is a path from j to i , and each link on the path was assigned a weight, the
20 function then returns the reciprocal of the sum of all links' weights on the path),
21 and Va is the importance vector for A to measure the rank of every object in A to
22 the set of core objects C .

23 Fig. 2 is a flowchart illustrating an example process 140 for community
24 mining. Process 140 is implemented by, for example, community mining module
25

1 104 of Fig. 1, and may be performed in software, hardware, firmware, or
2 combinations thereof.

3 Initially, one or more core sets of objects are identified (act 142). Each
4 core set of objects typically includes two or more objects, although alternatively a
5 core set may include a single object. The core sets are identified by identifying
6 groupings of objects with each object in a grouping being referenced by at least a
7 threshold number of other objects in the collection of objects. Once the core sets
8 of objects are identified, communities are created with the identified core sets (act
9 144). Each core set of objects identified in act 142 serves as the core or center of a
10 community.

11 Each community is then expanded, based on the core set of objects of the
12 community, by adding affiliated objects (act 146). Affiliated objects are objects
13 that have a link to one or more of the core set of objects in the community. These
14 affiliated objects may optionally be ranked in terms of importance (e.g., how well
15 each is deemed to relate to the community), as discussed in more detail below.

16 The communities created by identifying core sets in act 142 and expanding
17 with affiliated objects in act 146 can be further modified by performing core set
18 merging and/or community merging (act 148). This merging is optional. Core set
19 merging allows communities to be merged based on the similarity or overlap of
20 the core objects in the communities, while community merging allows
21 communities to be merged based on the similarity or overlap of all of the objects
22 in the communities. Core set merging and community merging are both discussed
23 in additional detail below.

24 Once the communities are created and expanded, and optionally merged,
25 the resulting communities are output as the one or more communities mined from

1 the collection of objects (act 150). Additionally, it should be noted that under
2 certain circumstances it is possible that the objects and links between objects are
3 such that no communities can be mined from the collection of objects.

4 Fig. 3 illustrates an example concentric ring model 180 of a community
5 mined from a collection of objects. The model 180 includes multiple concentric
6 rings 182, 184, 186, and 188. The center ring 182 includes the core objects 192 of
7 the community. The other objects illustrated in model 180 are affiliated objects
8 194, located in the various rings 184, 186, and 188 that are around the center ring
9 182. Any number of rings can be included in the concentric ring model 180 (e.g.,
10 as indicated by the ellipses between rings 186 and 188).

11 As can be readily seen from concentric ring model 180, the objects that are
12 deemed to be most important for the community are located in the center ring 182.
13 Other objects that are part of the community but that are deemed to be less
14 important are located in the various concentric rings 184, 186, and 188
15 surrounding center ring 182, with the objects that are located in rings closer to
16 center ring 182 deemed as being more important than rings located further from
17 center ring 182. Objects located in the same ring have the same importance level
18 to the community. Although the precise location of objects within the ring may
19 reveal a tiny variance in their importance to the community, they are deemed to be
20 the same from a macroscopical viewpoint.

21 In Fig. 3, the concentric rings are illustrated as circles. However, it should
22 be noted that the concentric ring model can be made up of concentric portions of
23 other geometric shapes as well (e.g., elliptical shapes, triangles, rectangles,
24 pentagons, etc.). Additionally, it should be noted that although the rings are
25 referred to herein as being concentric, the various rings may have the same center

1 or approximately the same center (that is, the rings need not have exactly the same
2 center).

3 It should also be noted that, rather than viewing the community as a
4 concentric ring model, other models may alternatively be used. For example, a
5 layered or stacked model may be used, with the core objects being at the bottom
6 (or top) of the stack and the affiliated objects being in higher (or lower) layers of
7 the stack.

8 Fig. 4 is a flowchart illustrating an example process 220 for finding core
9 sets of objects. Process 220 is implemented by, for example, core set
10 identification module 112 of Fig. 1, and may be performed in software, hardware,
11 firmware, or combinations thereof. Process 220 illustrates an example of act 142
12 of Fig. 2.

13 Initially, objects in the collection of objects and the link topology of the
14 collection of objects are identified (act 222). The link topology refers to which
15 objects in the collection are linked to which other objects in the collection.
16 Groups of objects that satisfy a link threshold are then identified (act 224). The
17 link threshold represents a minimum number of other objects in the collection that
18 must each link to a particular object in order for that object to be part of the group.
19 Multiple objects which both link to or cite the same other object are also referred
20 to as being co-linked (or co-cited) to that other object. For example, if the objects
21 are documents and the links are cites, and if the link threshold is two, then the
22 document groups are generated such that each document in a particular group is
23 cited by at least the same two other documents in the collection.

24 Once the groups are found in act 224, the largest groups of objects that are
25 not subsets of another group are identified as the core sets (act 226). It should be

1 noted that different core sets of different sizes can be mined from the same
2 collection of objects.

3 Fig. 5 shows an example that illustrates process 220 of Fig. 4. In the
4 example of Fig. 5, the collection 250 of objects includes six objects (A, B, C, D, E,
5 and F). Typically a collection would include more objects, but Fig. 5 is kept at six
6 for ease of explanation. Further, for ease of explanation assume that each of the
7 objects represents a document, and that the arrows represents links that are cites
8 from one document to another. The direction of the arrow indicates that one
9 document cites another (e.g., document F includes a cite to document C). Thus, it
10 can be seen in Fig. 5 that in the document collection 250: document A does not
11 cite any other document in collection 250; document B does not cite any other
12 document in collection 250; document C cites documents A and B; document D
13 cites documents A, B, and C; document E cites document A; and document F cites
14 document C.

15 Additionally, assume that in the example of Fig. 5, the link threshold is
16 two. Thus, the groups of objects found in act 224 that satisfy the link threshold of
17 two would be: the group of document A; the group of document B; the group of
18 document C; and the group of documents A and B. Although document C is cited
19 by two other documents (documents D and F), both of these two other documents
20 do not cite document A (and thus no group of documents A and C can be formed),
21 nor do both of these two other documents cite document B (and thus no group of
22 documents B and C can be formed).

23 Following this example, the group of documents A and B would be a core
24 set but the group of document A would not be a core set and the group of
25 document B would not be a core set (the group of document A is a subset of the

group of documents A and B, and the group of document B is a subset of the group of documents A and B). The group of document C would also be a core set (assuming single-object core sets are permitted), as the group of document C is not a subset of the group of documents A and B.

Returning to Fig. 4, the finding of groups of objects in act 224, as well as the identifying of the largest groups in act 226, can be performed in a variety of different manners. In one example implementation, the process is performed by identifying multiple groups of objects that may be core sets, and then refining these multiple groups by searching for larger groups and pruning out subsets of the larger groups. For example, the process may be performed by starting with single-object groups that satisfy the link threshold. These single-object groups are then combined into two-object groups that satisfy the link threshold, and any single-object groups that are subsets of the two-object groups are removed. This process continues until the largest group(s) of objects is found that satisfies (satisfy) the link threshold. Table I below includes example pseudo code for carrying out this process of acts 224 and 226.

Table I

1:	<i>Generate 1-itemsets IS₁ with minimal support S</i>
2:	<i>k ← 2</i>
3:	while <i>k ≤ m do</i>
4:	<i>Generate k-itemsets IS_k using (k-1)-itemsets IS_(k-1) with S</i>
5:	<i>Prun IS_(k-1) using IS_k</i>
6:	<i>k ← k + 1</i>
7:	end
8:	<i>Put IS₁ to IS_m to itemsets set IS</i>

1 In the pseudo code of Table I, the groups of objects are referred to as
2 itemsets, the notation " k -itemsets" refers to groups including k objects, and the
3 minimal support S refers to the link threshold.

4 As illustrated by the pseudo code of Table I, in line 1 all of the groups with
5 a single object that satisfy the link threshold are identified. The variable k is then
6 incremented to the value of two in line 2, and then a while loop spanning lines 3
7 through 8 begins. In the while loop of lines 3 through 8, groups of k objects (IS_k)
8 are generated by using combinations of the previously generated groups (IS_{k-1})
9 in line 4. All possible combinations of k objects from the objects of the IS_{k-1}
10 groups that satisfy the link threshold become groups of k objects (IS_k). So,
11 initially with k set to the value of two, groups of two objects are generated by
12 using combinations of the previously generated groups with one object (generated
13 in line 1). The groups generated in line 4 must satisfy the link threshold.

14 After the new groups are generated in line 4, groups with $k-1$ objects are
15 pruned in line 5 so that any of the groups with $k-1$ objects that are subsets of one
16 of the groups with k objects are removed. For example, if a group with document
17 A existed (a 1-object group), and a group with document B existed (also a 1-object
18 group), and a new group is generated with documents A and B (a 2-object group),
19 both of the 1-object groups would be pruned (removed). However, if a group with
20 document C also existed (a 1-object group), then this group would not be pruned
21 because it is not a subset of the 2-object group of documents A and B. This
22 pruning is performed because groups with more objects are more desired than
23 groups with fewer objects.

24 After pruning, the value of k is incremented by one. This process continues
25 in the while loop of lines 3 through 8 until a value m is reached. This value m

1 represents the longest itemset (the largest size group that satisfies the link
2 threshold). Once a value of k is reached for which no groups can be generated
3 having k objects that satisfy the link threshold, then the value of m is found (the
4 value of m then becomes $k-1$).

5 In line 8, the groups remaining when the while loop is exited (once the
6 value of m is hit) become the core sets. This will include at least one group with m
7 objects as well as possibly one or more other groups with fewer than m objects.
8 These different sized groups result because, as seen in the pseudo code of Table I,
9 the process begins with groups having single objects, and groups are removed in
10 line 5 if they are subsets of a larger group, but otherwise they are not removed.

11 The value of the link threshold can vary by implementation. In one
12 implementation, the value of the link threshold is determined empirically. In
13 another implementation, an initial estimation of the link threshold is determined as
14 follows. Initially, a number of objects from the collection are selected (e.g.,
15 randomly or pseudo randomly) to form the objects set R . The number selected can
16 vary, and in one example should be at least 1% of the total number of objects in
17 the collection. The number of objects linked to by each of these selected objects is
18 then identified, and the number of objects that link to each of these selected
19 objects is also identified. The amplified average links of each node can then be
20 used as follows to calculate the value for S (the link threshold):

$$21 \quad S = \frac{f \times \sum_R w_i}{\|R\|}$$

22
23 where f represents the amplifying frequency factor (e.g., set to 2 experimentally),
24 $\|R\|$ is the number of selected objects from the collection, and $\sum_R w_i$ is the weight
25

1 sum of all links related to R (that is, for any link in the graph, if there is a certain
2 object in R it connects to, then the weight on the link should be added to the sum).

3 Fig. 6 is a flowchart illustrating an example process 270 for expanding
4 communities based on core sets of objects. Process 270 is implemented by, for
5 example, community expansion module 114 of Fig. 1, and may be performed in
6 software, hardware, firmware, or combinations thereof. Process 270 illustrates an
7 example of act 146 of Fig. 2.

8 Initially, for a given core set of objects, all other objects in the collection of
9 objects (that is, all other objects in the collection of objects other than the given
10 core set of objects) that link to at least one object in the core set are identified as
11 an affiliated object (act 272). The community having that core set of objects is
12 then expanded to include the core set objects as well as the affiliated objects (act
13 274).

14 The affiliated objects are also ranked (act 276). The ranking of a particular
15 affiliated object is determined based on the number of objects in the core set that
16 the affiliated object links to – the larger the number of objects in the core set that
17 the affiliated object links to the higher its ranking is. For example, the affiliated
18 objects that link to all of the core objects may be given a rank of first, the affiliated
19 objects that link to one less than all of the core objects may be given a rank of
20 second, the affiliated objects that link to two less than all of the core objects may
21 be given a rank of third, and so forth. The ranking criteria for affiliated objects
22 can vary, as long as it could be used for sorting the affiliated objects and forming
23 the outer concentric rings 184, 186, and 188.

24 The affiliated objects are then assigned to particular ones of the concentric
25 rings based on their rankings (act 278). Affiliated objects with higher rankings are

1 assigned to rings closer to the center ring (where the core is located) than those
2 affiliated objects with lower rankings. For example, returning to Fig. 3, affiliated
3 objects with a rank of 1 may be assigned to ring 184, affiliated objects with a rank
4 of 2 may be assigned to ring 186, and so forth.

5 The ranking of affiliated objects can also be dependent on the weights of
6 the links between the affiliated objects and the objects in the core set. For
7 example, affiliated objects having higher-weighted links to the objects in the core
8 set may be given higher rankings than affiliated objects having lower-weighted
9 links. These link weights may be used to determine the rankings of the affiliated
10 objects, and/or to determine locations of objects within the concentric rings (e.g.,
11 affiliated objects having higher-weighted links to the objects in the core set are
12 located closer to the center ring (where the core is located) than affiliated objects
13 having lower-weighted links).

14 Thus, once the communities are created and expanded, the objects in the
15 communities that are deemed most important can be quickly and easily identified.
16 The most important or core objects are those in the center ring (the core set of
17 objects). With regard to the affiliated objects, the importance of the various
18 affiliated objects can be readily identified based on how close they are to the
19 center ring.

20 Fig. 7 is a flowchart illustrating an example process 300 for performing
21 core set merging. Process 300 is implemented by, for example, core set merging
22 module 116 of Fig. 1, and may be performed in software, hardware, firmware, or
23 combinations thereof. Process 300 illustrates an example of act 148 of Fig. 2.

24 Initially, core sets of two communities in the collection of objects are
25 identified (act 302). A check is then made as to whether there is sufficient overlap

or similarity of the identified core sets (act 304). The check as to whether there is sufficient overlap or similarity of the identified core sets is basically a check to determine whether the two core sets are similar enough that they should be combined into a single core set. Two core sets overlap if there are objects that are included in both core sets. An example of this situation is illustrated in Fig. 8.

In Fig. 8 three core sets of objects are illustrated: core set 320 including objects A, B, C, D, and E; core set 322 including objects A, B, C, and F; and core set 324 including objects D, E, and F. Given the overlapping of these core sets 320, 322, and 324, it may very well be desirable to combine these three core sets 320, 322, and 324 to generate a single core set. Furthermore, even if core set 324 did not exist, it may still be desirable to combine the two core sets 320 and 322.

Returning to Fig. 7, one or more rules (or constraints) are used to determine whether there is sufficient overlap or similarity of two core sets to justify merging the two core sets. In one example implementation, the following three constraints are defined to determine whether two core sets can be merged:

$$(1) \frac{\text{Min}(\|S_i\|, \|S_j\|)}{\|S_i \cap S_j\|} < 2$$

$$(2) \exists T \subset \|S_i \cup S_j - (S_i \cap S_j)\|, \frac{\|S_i \cup S_j - (S_i \cap S_j)\|}{\|T\|} < 2, \text{Support}(T) \geq S$$

$$(3) \|T\| \geq 2, \exists o_1 \in T \text{ and } o_1 \in (S_i - (S_i \cap S_j)), \exists o_2 \in T \text{ and } o_2 \in (S_j - (S_i \cap S_j))$$

where S_i represents the object set of a core set i , $\|S_i\|$ represents the number of objects in S_i , S_j represents the object set of a core set j , $\|S_j\|$ represents the number of objects in S_j , the *Min* operation returns the smallest of the values input to the *Min* operation (e.g., the smallest of $\|S_i\|$ and $\|S_j\|$), and *Support*(T) represents the support value of object set T (that is, the largest link threshold that T would

1 satisfy). T is a common subset of both S_i and S_j , and the calculated support value
2 of T should also meet the minimal support threshold S (e.g., as referenced above in
3 the pseudo code of Table I). If all three of these constraints are satisfied, then the
4 core set i and the core set j can be merged.

5 If the two identified core sets can be merged, then the two communities
6 having those two core sets are merged, resulting in a single community (act 306).
7 All of the affiliated objects in the communities of each of the two identified core
8 sets become affiliated objects in the new single community (unless one of the
9 affiliated objects becomes a core object). The rankings for the affiliated objects (if
10 any) may optionally also be re-determined in act 306.

11 A check is then made as to whether there are any additional core sets to
12 check for merging (act 308). The check is also made if the two identified core sets
13 cannot be merged (from act 304). In one implementation, process 300 checks all
14 combinations of two core sets to determine whether any of the combinations can
15 be merged. When a new community is generated by core set merging, the core set
16 of this new community may also be used as one of the two core sets when
17 checking all of these combinations. If there are additional combinations of core
18 sets to check, then process 300 returns to act 302 where two more core sets are
19 identified. However, if there are no more combinations of core sets to check, then
20 the core set merging is finished (act 310).

21 It should be noted that, when two core sets are merged using process 300,
22 the link threshold discussed above is no longer satisfied by the merged core set (if
23 it were satisfied, then the merged core set should have been identified in the
24 processes for finding core sets discussed above).

Fig. 9 is a flowchart illustrating an example process 350 for performing community merging. Process 350 is implemented by, for example, community merging module 118 of Fig. 1, and may be performed in software, hardware, firmware, or combinations thereof. Process 350 illustrates an example of act 148 of Fig. 2.

Initially, two communities in the collection of objects are identified (act 352). A check is then made as to whether there is sufficient overlap or similarity of the identified communities (act 354). The check as to whether there is sufficient overlap or similarity of the identified communities is basically a check to determine whether the two communities are similar enough that they should be combined into a single community, even though their core sets may be different. Two communities overlap if there are objects that are included in both communities. An example of this situation is illustrated in Fig. 10.

In Fig. 10 two communities are illustrated. The communities have different core sets, but do have some overlapping affiliate objects. The overlapping affiliate objects are illustrated in Fig. 10 as cross-hatched. Given the overlapping of these two communities, it may very well be desirable to combine the two communities into a single community.

Returning to Fig. 9, one or more rules (or constraints) are used to determine whether there is sufficient overlap or similarity of two communities to justify merging the two communities. In one example implementation, the following constraint is defined to determine whether two communities can be merged:

$$\frac{\text{Min}\left(\sum_{\|ESi\|} w_k, \sum_{\|ESj\|} w_k\right)}{\sum_{\|ESi \cap ESj\|} w_k} < 2$$

1
2 where ES_i represents the affiliated object set expanded from the core set S_i , ES_j
3 represents the affiliated object set expanded from the core set S_j , w_k represents the
4 rank of an affiliated object, and the *Min* operation returns the smallest of the
5 values input to the *Min* operation.

6 If the two identified communities can be merged, then the two communities
7 are merged, resulting in a single community (act 356). All of the affiliated objects
8 in the communities of each of the two identified core sets become affiliated
9 objects in the new single community (unless one of the affiliated objects becomes
10 a core object). The rankings for the affiliated objects (if any) may optionally also
11 be re-determined in act 356.

12 A check is then made as to whether there are any additional communities to
13 check for merging (act 358). The check is also made if the two identified
14 communities cannot be merged (from act 354). In one implementation, process
15 350 checks all combinations of two communities to determine whether any of the
16 combinations can be merged. When a new community is generated by community
17 merging, this new community may also be used as one of the two communities
18 when checking all of these combinations. If there are additional combinations of
19 communities to check, then process 350 returns to act 352 where two more
20 communities are identified. However, if there are no more combinations of
21 communities to check, then the community merging is finished (act 360).

22 It should be noted that, analogous to the core set merging discussed above,
23 when two communities are merged using process 350, the link threshold discussed
24 above is no longer satisfied by the core set of the merged community (if it were
25

1 satisfied, then the merged community should have been identified in the processes
2 for finding core sets discussed above).

3 It should also be noted that, as can be seen from the description herein,
4 there is no limit as to the number of different communities an object can belong to.
5 For example, an object may be an affiliate object in multiple communities, an
6 object may be an affiliate object in one or more communities and a core object in
7 one or more other communities, an object may be a core object in multiple
8 communities, and so forth.

9 It should further be noted that, rather than identifying large groups of
10 objects during core set identification (e.g., as discussed above with respect to
11 Fig. 4 and the pseudo code of Table I), small groups of objects may alternatively
12 be identified. For example, the core set identification may simply identify groups
13 with two or three objects as core sets, without attempting to find groups with
14 larger numbers of objects. After these smaller groups are identified as core sets,
15 the core set merging of Fig. 7 and community merging of Fig. 9 can be relied on to
16 merge the communities.

1 As can be seen from the description herein, the community mining based
2 on core objects and affiliated objects described herein can have several
3 characteristics. Some of these characteristics are as follows:

- 4 • Core objects and affiliated objects in a community are distinguished.
5 This allows the objects deemed as being most representative of the
6 community (the core objects) to be highlighted and further allows the
7 affiliated objects to be ranked according to their deemed importance to
8 the core objects.
- 9 • The core of a community is made up of one or more objects. In many
10 situations, the true core of a community is often a combination of
11 multiple objects. By allowing the core to be made up of multiple
12 objects, more coherent communities can be created.
- 13 • The objects in the core of a community are not required to be tightly
14 linked (there is no requirement as to direct links among the objects in
15 the core). In fact, it is possible for none of the objects in the core set of
16 a community to directly link to other objects in the core set of the
17 community.
- 18 • Objects are part of a core set of a community based on the links to those
19 objects, not based on how many other objects they may link to.
- 20 • Each affiliated object is ranked according to how many of the core
21 objects of the community the affiliated object is linked to. The more
22 core objects in a community an affiliated object links to, the better it is
23 deemed to match the topic of the community.

24 Fig. 11 illustrates a general computer environment 400, which can be used
25 to implement the techniques described herein. The computer environment 400 is

1 only one example of a computing environment and is not intended to suggest any
2 limitation as to the scope of use or functionality of the computer and network
3 architectures. Neither should the computer environment 400 be interpreted as
4 having any dependency or requirement relating to any one or combination of
5 components illustrated in the exemplary computer environment 400.

6 Computer environment 400 includes a general-purpose computing device in
7 the form of a computer 402. Computer 402 can implement, for example,
8 community mining module 104 of Fig. 1. The components of computer 402 can
9 include, but are not limited to, one or more processors or processing units 404, a
10 system memory 406, and a system bus 408 that couples various system
11 components including the processor 404 to the system memory 406.

12 The system bus 408 represents one or more of any of several types of bus
13 structures, including a memory bus or memory controller, a peripheral bus, an
14 accelerated graphics port, and a processor or local bus using any of a variety of
15 bus architectures. By way of example, such architectures can include an Industry
16 Standard Architecture (ISA) bus, a Micro Channel Architecture (MCA) bus, an
17 Enhanced ISA (EISA) bus, a Video Electronics Standards Association (VESA)
18 local bus, and a Peripheral Component Interconnects (PCI) bus also known as a
19 Mezzanine bus.

20 Computer 402 typically includes a variety of computer readable media.
21 Such media can be any available media that is accessible by computer 402 and
22 includes both volatile and non-volatile media, removable and non-removable
23 media.

24 The system memory 406 includes computer readable media in the form of
25 volatile memory, such as random access memory (RAM) 410, and/or non-volatile

1 memory, such as read only memory (ROM) 412. A basic input/output system
2 (BIOS) 414, containing the basic routines that help to transfer information
3 between elements within computer 402, such as during start-up, is stored in ROM
4 412. RAM 410 typically contains data and/or program modules that are
5 immediately accessible to and/or presently operated on by the processing unit 404.

6 Computer 402 may also include other removable/non-removable,
7 volatile/non-volatile computer storage media. By way of example, Fig. 11
8 illustrates a hard disk drive 416 for reading from and writing to a non-removable,
9 non-volatile magnetic media (not shown), a magnetic disk drive 418 for reading
10 from and writing to a removable, non-volatile magnetic disk 420 (e.g., a “floppy
11 disk”), and an optical disk drive 422 for reading from and/or writing to a
12 removable, non-volatile optical disk 424 such as a CD-ROM, DVD-ROM, or other
13 optical media. The hard disk drive 416, magnetic disk drive 418, and optical disk
14 drive 422 are each connected to the system bus 408 by one or more data media
15 interfaces 426. Alternatively, the hard disk drive 416, magnetic disk drive 418,
16 and optical disk drive 422 can be connected to the system bus 408 by one or more
17 interfaces (not shown).

18 The disk drives and their associated computer-readable media provide non-
19 volatile storage of computer readable instructions, data structures, program
20 modules, and other data for computer 402. Although the example illustrates a hard
21 disk 416, a removable magnetic disk 420, and a removable optical disk 424, it is to
22 be appreciated that other types of computer readable media which can store data
23 that is accessible by a computer, such as magnetic cassettes or other magnetic
24 storage devices, flash memory cards, CD-ROM, digital versatile disks (DVD) or
25 other optical storage, random access memories (RAM), read only memories

1 (ROM), electrically erasable programmable read-only memory (EEPROM), and
2 the like, can also be utilized to implement the exemplary computing system and
3 environment.

4 Any number of program modules can be stored on the hard disk 416,
5 magnetic disk 420, optical disk 424, ROM 412, and/or RAM 410, including by
6 way of example, an operating system 426, one or more application programs 428,
7 other program modules 430, and program data 432. Each of such operating
8 system 426, one or more application programs 428, other program modules 430,
9 and program data 432 (or some combination thereof) may implement all or part of
10 the resident components that support the distributed file system.

11 A user can enter commands and information into computer 402 via input
12 devices such as a keyboard 434 and a pointing device 436 (e.g., a “mouse”).
13 Other input devices 438 (not shown specifically) may include a microphone,
14 joystick, game pad, satellite dish, serial port, scanner, and/or the like. These and
15 other input devices are connected to the processing unit 404 via input/output
16 interfaces 440 that are coupled to the system bus 408, but may be connected by
17 other interface and bus structures, such as a parallel port, game port, or a universal
18 serial bus (USB).

19 A monitor 442 or other type of display device can also be connected to the
20 system bus 408 via an interface, such as a video adapter 444. In addition to the
21 monitor 442, other output peripheral devices can include components such as
22 speakers (not shown) and a printer 446 which can be connected to computer 402
23 via the input/output interfaces 440.

24 Computer 402 can operate in a networked environment using logical
25 connections to one or more remote computers, such as a remote computing device

1 448. By way of example, the remote computing device 448 can be a personal
2 computer, portable computer, a server, a router, a network computer, a peer device
3 or other common network node, and the like. The remote computing device 448 is
4 illustrated as a portable computer that can include many or all of the elements and
5 features described herein relative to computer 402.

6 Logical connections between computer 402 and the remote computer 448
7 are depicted as a local area network (LAN) 450 and a general wide area network
8 (WAN) 452. Such networking environments are commonplace in offices,
9 enterprise-wide computer networks, intranets, and the Internet.

10 When implemented in a LAN networking environment, the computer 402 is
11 connected to a local network 450 via a network interface or adapter 454. When
12 implemented in a WAN networking environment, the computer 402 typically
13 includes a modem 456 or other means for establishing communications over the
14 wide network 452. The modem 456, which can be internal or external to computer
15 402, can be connected to the system bus 408 via the input/output interfaces 440 or
16 other appropriate mechanisms. It is to be appreciated that the illustrated network
17 connections are exemplary and that other means of establishing communication
18 link(s) between the computers 402 and 448 can be employed.

19 In a networked environment, such as that illustrated with computing
20 environment 400, program modules depicted relative to the computer 402, or
21 portions thereof, may be stored in a remote memory storage device. By way of
22 example, remote application programs 458 reside on a memory device of remote
23 computer 448. For purposes of illustration, application programs and other
24 executable program components such as the operating system are illustrated herein
25 as discrete blocks, although it is recognized that such programs and components

1 reside at various times in different storage components of the computing device
2 402, and are executed by the data processor(s) of the computer.

3 Various modules and techniques may be described herein in the general
4 context of computer-executable instructions, such as program modules, executed
5 by one or more computers or other devices. Generally, program modules include
6 routines, programs, objects, components, data structures, etc. that perform
7 particular tasks or implement particular abstract data types. Typically, the
8 functionality of the program modules may be combined or distributed as desired in
9 various embodiments.

10 An implementation of these modules and techniques may be stored on or
11 transmitted across some form of computer readable media. Computer readable
12 media can be any available media that can be accessed by a computer. By way of
13 example, and not limitation, computer readable media may comprise “computer
14 storage media” and “communications media.”

15 “Computer storage media” includes volatile and non-volatile, removable
16 and non-removable media implemented in any method or technology for storage
17 of information such as computer readable instructions, data structures, program
18 modules, or other data. Computer storage media includes, but is not limited to,
19 RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM,
20 digital versatile disks (DVD) or other optical storage, magnetic cassettes, magnetic
21 tape, magnetic disk storage or other magnetic storage devices, or any other
22 medium which can be used to store the desired information and which can be
23 accessed by a computer.

24 “Communication media” typically embodies computer readable
25 instructions, data structures, program modules, or other data in a modulated data

1 signal, such as carrier wave or other transport mechanism. Communication media
2 also includes any information delivery media. The term “modulated data signal”
3 means a signal that has one or more of its characteristics set or changed in such a
4 manner as to encode information in the signal. By way of example, and not
5 limitation, communication media includes wired media such as a wired network or
6 direct-wired connection, and wireless media such as acoustic, RF, infrared, and
7 other wireless media. Combinations of any of the above are also included within
8 the scope of computer readable media.

9 Various flowcharts are described herein and illustrated in the
10 accompanying Figures. The ordering of acts in these flowcharts are examples
11 only – these orderings can be changed so that the acts are performed in different
12 orders and/or concurrently.

13 Although the description above uses language that is specific to structural
14 features and/or methodological acts, it is to be understood that the invention
15 defined in the appended claims is not limited to the specific features or acts
16 described. Rather, the specific features and acts are disclosed as exemplary forms
17 of implementing the invention.